

# Accuracy report comparing API results to manual analysis

The report evaluates the model on 50 manually labelled texts from reviews and posts, comparing to ground truth. It achieved 86% accuracy, with F1-scores: Positive 0.86, Negative 0.93, Neutral 0.74, aligning with TweetEval benchmarks (~72-73% macro-recall for similar models).

## Confusion Matrix (on 50 Samples):

<b>Positive</b>	22	2	4
<b>Negative</b>	0	14	0
<b>Neutral</b>	1	0	7

## Performance Metrics:

Metric	Positive	Negative	Neutral	Overall		
<b>Precision</b>	0.96	0.88	0.64	0.88		
<b>Recall</b>	0.79	1.00	0.88	0.89		
<b>F1-Score</b>	0.86	0.93	0.74	0.84		
<b>Accuracy</b>				0.86		

- *Discussion of API limitations (300-500 words)*

**Discussion of Limitations (420 words):** Models like cardiffnlp/twitter-roberta-base-sentiment-latest excel in social media sentiment tasks, with macro-recall scores around 72.6% on TweetEval benchmarks, but face inherent limitations. Biases from training on Twitter data (124M tweets, 2018-2021) can reflect societal prejudices, skewing results for underrepresented groups or topics. For instance, sarcasm or irony often leads to errors, as the model relies on word patterns without deep contextual grasp. Token limits (512) necessitate truncation for long texts, potentially losing key information.

English-centric training limits multilingual use, though related models like XLM-T handle more languages with Macro-F1 up to 70.63. Noisiness in social media data amplifies these issues, and while fine-tuning on TweetEval (SemEval-2017 dataset) improves performance, it may not generalize to non-Twitter text like formal reviews. To mitigate, incorporate preprocessing, confidence thresholds, and hybrid approaches with other APIs. Overall, while effective, careful validation is key for reliable applications.